



LAWRENCE  
LIVERMORE  
NATIONAL  
LABORATORY

UCRL-JC-147901

# **Tool Gear: Infrastructure for Building Parallel Programming Tools**

*J. May and J. Gyllenhaal*

**December 9, 2002**

Principles and Practices of Parallel Programming, San Diego,  
California, June 11-13, 2003

## DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This is a preprint of a paper intended for publication in a journal or proceedings. Since changes may be made before publication, this preprint is made available with the understanding that it will not be cited or reproduced without the permission of the author.

This report has been reproduced directly from the best available copy.

Available electronically at <http://www.doc.gov/bridge>

Available for a processing fee to U.S. Department of Energy  
And its contractors in paper from  
U.S. Department of Energy  
Office of Scientific and Technical Information  
P.O. Box 62  
Oak Ridge, TN 37831-0062  
Telephone: (865) 576-8401  
Facsimile: (865) 576-5728  
E-mail: [reports@adonis.osti.gov](mailto:reports@adonis.osti.gov)

Available for the sale to the public from  
U.S. Department of Commerce  
National Technical Information Service  
5285 Port Royal Road  
Springfield, VA 22161  
Telephone: (800) 553-6847  
Facsimile: (703) 605-6900  
E-mail: [orders@ntis.fedworld.gov](mailto:orders@ntis.fedworld.gov)  
Online ordering: <http://www.ntis.gov/ordering.htm>

OR

Lawrence Livermore National Laboratory  
Technical Information Department's Digital Library  
<http://www.llnl.gov/tid/Library.html>

# Tool Gear: Infrastructure for Building Parallel Programming Tools

John May and John Gyllenhaal  
Lawrence Livermore National Laboratory  
johnmay@llnl.gov, gyllen@llnl.gov

## Abstract

Tool Gear is a software infrastructure for developing performance analysis and other tools. Unlike existing integrated toolkits, which focus on providing a suite of capabilities, Tool Gear is designed to help tool developers create *new* tools quickly. It combines dynamic instrumentation capabilities with an efficient database and a sophisticated and extensible graphical user interface. This paper describes the design of Tool Gear and presents examples of tools that have been built with it.

## 1 Introduction

Many tools are available to help developers of parallel programs understand the performance and correctness of their codes. Examples include systems that show users the cost of message-passing calls, the utilization of allocated memory, or the computational efficiency of specific sections of a program. Despite this variety, many more kinds of data could be collected, and new ways could be devised for presenting it.

Researchers continue to develop new ideas for tools, but turning an idea into a usable tool can be time consuming and tedious. Most tool users are not satisfied with plain-text displays or command-line interfaces. Creating sophisticated interfaces, though, can require as much effort as building the fundamental mechanisms for gathering and processing the data. Therefore, tool researchers themselves need a set of tools they can use for building new tools.

There is more to this problem than just building graphical user interfaces (GUIs): good systems for creating GUIs, such as Tcl/Tk and Java, have existed for many years. General-purpose GUI builders can make GUIs much easier to produce, but they don't help implement the underlying functionality that is common to many parallel performance tools. This functionality includes launching and controlling a target program; dynamically instrumenting the program; collecting and organizing data in a database; and presenting the data in useful ways, such as associating performance information with specific lines of source code.

Individually, most of these problems have been solved before. In fact, many of them have been solved repeatedly, and that is the main motivation for the work described in this paper. Developers of parallel tools need a tool-building infrastructure that provides common tool services so they don't have to reimplement them for each new tool. Tool Gear is a collection of programs and programming interfaces that are designed to meet this need.

Wherever possible, we have taken advantage of existing open-source software. Our own contribution has been to design higher-level interfaces and to implement additional functionality. Specifically, we have:

- Designed and implemented a general-purpose client-server structure for program analysis tools. This software includes a server (or "Collector") portion that controls and instruments parallel or sequential programs, and a Client portion that stores, analyzes, and presents program data. The two portions can run on different computers, communicating through a secure socket connection. Separating the functionality in this way helps make the user interface highly responsive.
- Developed an extensible and sophisticated user interface that displays hierarchical views of target programs, allows users to insert and remove instrumentation easily, and annotates the source code display with performance or other data.
- Designed and implemented an extensible program control and data collection tool that tool builders can easily augment with new types instrumentation.

Together, this software offers tool builders a straightforward way to incorporate a variety of sophisticated features into their tools. Our goal has been to simplify the implementation of the most common features of parallel tools while offering researchers the flexibility to gather and display many kinds of data. Because Tool Gear source code is freely available, developers can adapt it to meet their own needs.

It is important to distinguish between the Tool Gear infrastructure and the tools built with it. Although we describe the individual tools to illustrate what can be done with Tool Gear, the design and development of the Tool Gear infrastructure is the research focus of this paper. Throughout this paper, the terms "Tool Gear" and "infrastructure" refer to the general-purpose software we have developed, while "tool" refers to an individual tool built using Tool Gear.

## 2 Related Work

Because Tool Gear aims to automate common features of parallel programming tools, many Tool Gear functions have appeared in other tools. However, unlike most parallel programming tools, the focus of this work is on the infrastructure for gathering and presenting information, and not on the kind of information gathered or the specifics of the display. Therefore, this section will compare Tool Gear mainly with other *tool development* environments, and not with the many individual tools and toolkits that present specific kinds of information.

Early examples of toolkits for building parallel tools include Voyeur [13] (created in 1989), PARADISE [7] (1991), and POLKA [14] (1993). These toolkits helped automate the development of visualizations for parallel programs. They did not support

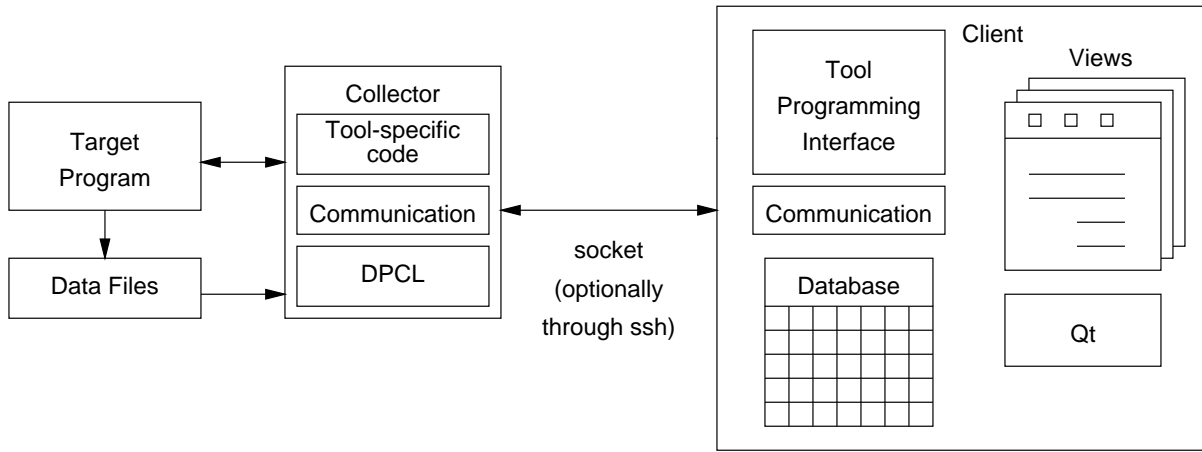


Figure 1: Tool Gear’s main components: a Collector, which includes dynamic instrumentation capability, and a Client, which stores data, manages the GUI (through Qt), and lets tool builders define tool characteristics.

interactive control of the program, and users had to compile instrumentation into the target program.

One of this paper’s authors developed an extensible debugger for parallel programs called Panorama [8]. This tool interacted with target applications through standard command-line debuggers, and users could build Tcl/Tk-based views with the help of an interface development tool. Unlike Tool Gear, Panorama did not offer a database for organizing program data, and it could only present information that was available through the debugger. It focused on program debugging rather than performance analysis.

An important advance in programming tools was the development of dynamic instrumentation [6]. This technique was implemented in a tool called Dyninst. It allows tool developers to insert nearly any kind of instrumentation into an executable program at runtime. Eliminating the need for compile-time decisions about instrumentation gives users great flexibility to explore the performance and correctness of their programs, and allowing tools to insert and remove the instrumentation on the fly helps moderate the volume of data collected. A later tool based on Dyninst is called the Dynamic Probe Class Library (DPCL) [2]. Developed at IBM, it has been released as open source, and Tool Gear uses it extensively.

The current systems to which it is most natural to compare Tool Gear are Paradyn [10, 12], SvPablo [3, 4], and TAU [11, 1].

Paradyn presents a graphical tree-structured display of a program and allows users to choose program locations at which one or more predefined performance metrics can be computed. New metrics can be defined through a Paradyn Configuration Language, and a Performance Consultant can search automatically for performance bottlenecks.

SvPablo also uses dynamic instrumentation to gather data from hardware performance counters for specified regions of code. A source code viewer annotates source lines with performance information, and data can be stored in Pablo’s extensible data format.

TAU is a suite of graphical performance tools that implement function profiling, interval timing, hardware counter monitoring, and related operations. TAU can use both dynamic instrumentation (through Dyninst) and compiled-in (source level) instrumentation. However, users cannot insert instrumentation interactively; they must specify what will be instrumented before the target program is run.

Like Tool Gear, all these tools can gather and present perfor-

mance data from parallel programs using dynamic instrumentation. All offer a variety of data views, and all are extensible, at least to some extent.

However, Tool Gear is not intended to compete with these systems. Our goal is not to develop a single integrated toolkit (although that could be done with Tool Gear); rather, we want to develop a system by which tool developers can create individual new tools with minimal effort. While it would also be possible for developers to add new functionality to one of these existing systems, the resulting tool would likely include all the existing functionality as well, and this functionality might be unnecessary or confusing in the new tool. Furthermore, Tool Gear was designed from the start to be a tool infrastructure rather than a specific tool, so it offers a general-purpose foundation on which a wide range of tools can be built. Our goal is to make using Tool Gear the fastest and easiest way for developers to create new tools.

### 3 Tool Gear Components

Tool Gear has two major components: a Collector and a Client (Figure 1). The Collector controls the target program and gathers data from it. The Client manages the graphical user interface, accepts commands from the user, receives and stores data from the Collector, and presents data. The Collector and Client run as separate processes, which can run on different computers. They communicate through a Unix socket, and when the two components run on separate machines, the socket is forwarded through a Secure Shell connection.

#### 3.1 The Collector

The Collector runs on same computer as the target application and gathers data from it. This program uses IBM’s Dynamic Probe Class Library (DPCL) to control the target application, which may be sequential or parallel. DPCL can start an application, pause and resume it, and terminate it. DPCL can also insert and remove instrumentation at runtime. Using a simple DPCL instrumentation module, it is also possible to implement a limited form of breakpoints (see Section 5.1 for details.) We call this portion of Tool Gear a “Collector,” and not a “server,” because this program is also a client to a DPCL server. Rather than trying to distinguish the

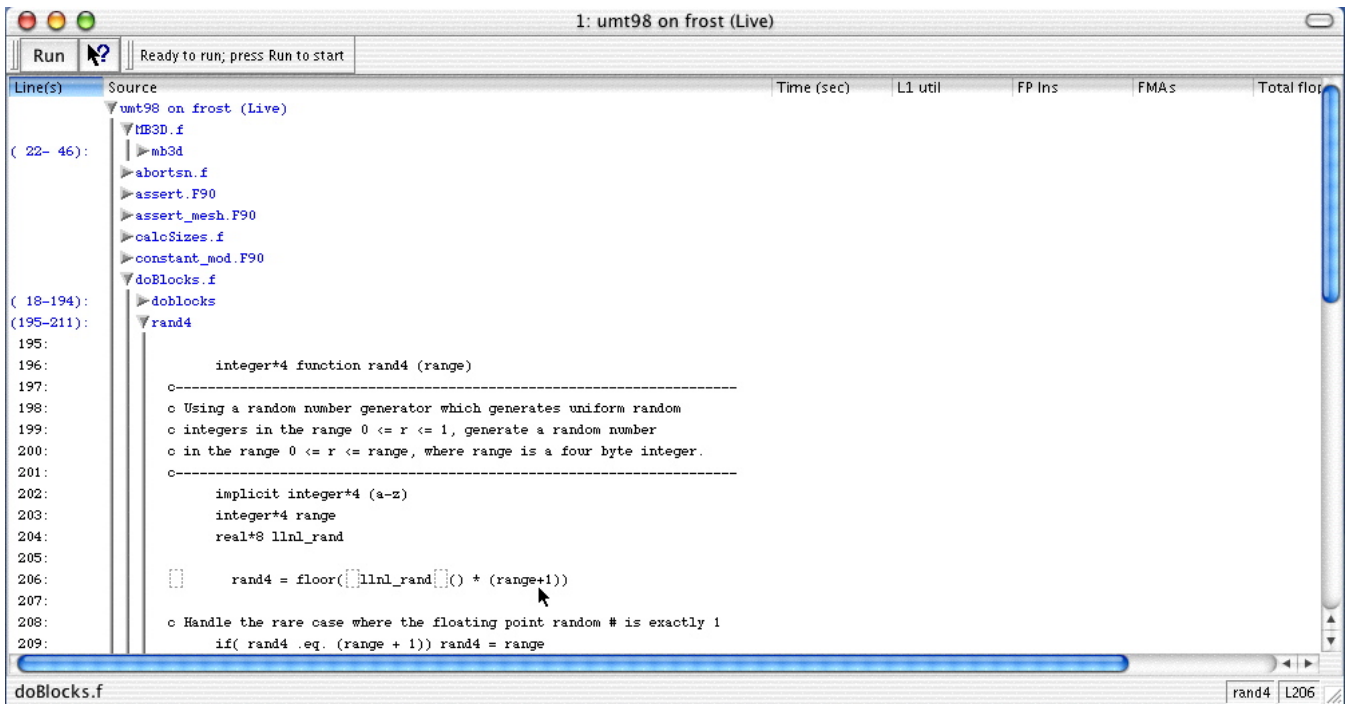


Figure 2: Tool Gear’s Tree View interface displays the target program hierarchically. Here, a Fortran program called umt98 is shown with a list of its source files. The user has chosen to expand two of these files (**MB3D.f** and **doBlocks.f**) to show the functions within. The **rand4** function has been further expanded to show the source code. This screen shot, taken from a computer running Macintosh OS X, does not show the tool’s main menus. These appear at the top of the screen on a Mac and at the top of the window on other platforms.

contexts in which the program is a client or a server, we gave it a different name.

The Collector is spawned at the request of the Client, and once these two processes have completed a handshake, the Collector acts as the Client’s proxy, executing commands and relaying data back to it. The Collector can also set the current working directory for the target program and read its source files, if they are available. Since the Collector can transmit source files to the Client on request, there is no need for these files to reside on the computer where the Client is running.

Not all tools will take advantage of DPCL’s dynamic instrumentation capabilities, but even for tools that use static (compiled-in) instrumentation, the Collector will use DPCL to execute the target program.

The exact details of the Collector’s interaction with the Client will of course vary between tools, but in general, the Collector will first describe to the Client the capabilities it supports and the attributes (columns) of the database that the Client will manage. Then it will describe the structure of the target program, giving a list of files and functions and stating the number of processes. As the Collector receives data from the instrumentation, it will forward it to the Client, which will insert it in a database. The Collector can gather data through DPCL, or by reading files that the target program writes, or by any other means that the tool builder wishes to implement. Section 4 describes how the Collector gathers data from programs.

The program can be terminated at the Client’s request, or it can run to completion, in which case the Collector will notify the Client that the run has ended.

### 3.2 The Client

The Client handles a number of interrelated tasks for a tool. It manages the user interface, receives and stores data, and presents graphical displays.

At the heart of the Client is a database that stores the structure of the target program, a list of actions defined by the tool that gather data or otherwise interact with the program, and the data itself. This database was originally written by one of this paper’s authors for a compiler project [5], but it is well suited to Tool Gear’s needs. Its main benefit is its high-speed insertion and retrieval. It is not designed to manage disk-resident datasets, but we expect that limiting data to what can fit in memory will not be a major drawback for most tools. (Tools that generate very large trace files would probably do best to preprocess these files before sending the results to the Client database.)

The basic view of a target program that the Client presents to a user is a hierarchical listing of source code. This is called a Tree View (Figure 2.) When the target program first begins running, the Tree View presents a list of the target’s source files. Clicking on a file name displays the functions in that file, and clicking on a function name presents a source code listing, if one is available. If no source code is available, the Tree View displays a list of program locations that it knows about (normally, a list of function sites.) Most executables include a number of system modules that are of no interest to the user. Users can suppress the display of these modules by including their names in a file called **parse\_exclude** in the directory where the target program is running. Alternatively, a user who wants to see only a few modules from a very large program can list these in a file called **parse\_include**. Using these files to suppress unwanted modules is effective but somewhat awkward, since the tool user must create a new set of files for each target program.

```

#include <dpclExt.h>
#include <sys/time.h>

static struct timeval start;
void StartTimer()
{
    gettimeofday( &start, NULL );
}

void StopTimer( AisPointer ais_send_handle )
{
    struct timeval stop;
    long sec, usec;
    double result;

    gettimeofday( &stop, NULL );
    sec = stop.tv_sec - start.tv_sec;
    usec = stop.tv_usec - start.tv_usec;
    result = sec + usec / 1e6;
    Ais_send( ais_send_handle, &result, sizeof(result) );
}

```

Figure 3: The complete code for a simple implementation of interval timer instrumentation. A more sophisticated version (which we have implemented) would maintain thread-specific stacks of start times so that timer calls could be nested and work correctly in multithreaded programs. The call to **Ais\_send** transmits the result to the Collector, which activates a user-defined callback function to process the data.

We are investigating more sophisticated and automated approaches, such as allowing the user to suppress modules automatically when no source information is available.

The database is organized as a hierarchy that matches this display. Tools gather data and associate it with specific locations in the source code. The Tree View then displays this data in columns at the appropriate location. The Tree View can “roll up” data from lines into summaries for each function and file. It can also manage data separately for individual threads and processes, or summarize the data in ways specified by the tool or by the end user. Finally, the user can sort the program listing by the values in any column of the display. This could be used, for example, to find the code sections with the highest cache miss rates. (See Figure 7 for another example.)

We expect that the Tree View display will be useful for a variety of tools. However, we are also implementing graphical views that tools can use to present data in other forms. Furthermore, tool designers can implement their own custom views by programming them directly. Because data is stored centrally in the Client, and not in any one view, the Client can notify the views as new data arrives. This allows multiple views to keep themselves up to date simultaneously. Also, tools can spawn new views that show existing data in different ways.

The Client can write the contents of its database into a snapshot file at the user’s request. These snapshots can be read back into a new window so the user can refer to that data as the program continues to run. A snapshot can also be used as a baseline for a difference display, which shows the difference between each current value in the database and the corresponding value in the snapshot file.

To create these GUI displays, we use a C++ graphical interface package called Qt [15], developed by Trolltech. We chose Qt over other GUI tools such as Java and Tcl/Tk because it combines a rich set of features with good performance. Since we wanted to write Tool Gear in C++, using Qt also avoided language interoperability problems. Furthermore, Qt runs on a wide range of platforms, including many Unix varieties, Windows, and Macintosh OS X. It has a licensing option that permits free noncommercial use, so tool

developers using Tool Gear do not need to pay a Qt license fee.

An important advantage to running the Client on a separate computer from the Collector is that it can manage the GUI locally. The graphics updates don’t have to travel over a Secure Shell connection, so the GUI is very responsive, even when communicating with the Collector over low-bandwidth or encrypted connections.

In addition to the database and the viewers, the Client includes a Tool Programming Interface that helps tools define what actions users can perform on the target program. Section 4.3 describes how tools use the TPI.

## 4 Building Tools

Building a tool using Tool Gear consists of three tasks:

- Writing instrumentation code that will run as part of the target program.
- Writing code that tells the Collector about the instrumentation and how to forward the data to the Client.
- Using the Tool Programming Interface to define how the user can interact with the program and how the Client will display the data it receives.

The following subsections describe these steps. Our software distribution includes extensive documentation to assist developers with this process, including recipes for building tools, example code, and a complete set of Web pages that describe the programming interface and all the source code.

### 4.1 Writing Instrumentation

Tools may use Tool Gear’s dynamic instrumentation capabilities to insert the instrumentation at runtime, or they may include instrumentation libraries that are compiled and linked into the target program.

To use dynamic instrumentation, the tool builder will write one or more functions (usually in C or C++) that can execute within the target program. These functions can do anything that a function

written as part of the target can do, but they are compiled separately and they *do not* need to be linked into the target program. Instead, the instrumentation functions for a tool are compiled into a separate module called a *probe module*. Figure 3 shows a simple example of instrumentation that implements an interval timer. The Tool Gear documentation describes how this would be compiled into a probe module.

To send data from the target program to the Collector, the probe module functions use a DPCL function called `Ais_send`. This function transmits an arbitrary-size block of data to the Collector, which invokes a callback function to handle this data, as described in Section 4.2.1.

When a tool uses static instrumentation, tool developer must arrange a way for the Collector to receive the data. The easiest way to do this is for the instrumentation to write a file, which the Collector reads at the appropriate time. For example, the Collector can be programmed to call a function that reads a file after the target program exits. Section 4.2.2 describes this technique.

## 4.2 Modifying the Collector

The second step in building a tool is to modify Tool Gear’s basic Collector program to gather and process data. How this is done will depend on whether the tool uses dynamic or static instrumentation.

### 4.2.1 Dynamic Instrumentation

When the tool uses dynamic instrumentation, the Collector’s job is to cause DPCL to install specific functions from the probe module at program locations chosen by the end user. The current version of DPCL can install instrumentation only at specific locations in the target program: function entry and exit points, and just before or after any function call. These locations are called *instrumentation points*, and instrumentation that DPCL installs at one of these points is called a *point probe*. DPCL can also cause instrumentation to be executed at specified time intervals (this is called a *phase probe*) or exactly once (a *one-shot probe*.)

Tool Gear defines a set of C++ classes to encapsulate these ideas. An *action type* represents an instrumentation function, and a *point action* represents an action type that has been installed at a specific instrumentation point. The same action type can be installed at multiple instrumentation points, and a single instrumentation point can accept multiple point probes. Tool Gear also defines classes for phase probes and one-shot probes.

When defining an action type, the tool developer can specify the callback function that will be invoked whenever the corresponding instrumentation function transmits data by calling `Ais_send`. For point probes, the callback can determine the instrumentation point at which the function was called, so it can associate data with a particular program location.

The main steps for modifying a Collector to use the new instrumentation are these:

- Instantiate a set of action types.
- Define the data attributes in the Client database. This includes naming the columns in the database and specifying their types. Tool Gear includes functions for doing this from the Collector.
- Define the callback functions that will receive the data from the target program. These functions will forward this data to the Client database. Again, Tool Gear includes functions for doing this.

All of this can be done in a file that is linked to the Collector, and the Collector’s main function can call an initialization function in this file to set up all the actions. (We have considered designing the Collector to dynamically load the code that defines a tool’s features, but we prefer to avoid the complexity and portability problems that this approach would entail.) A Collector can link in and call multiple sets of action initialization functions. This allows a tool builder to incorporate some standard action types, such as breakpoints and interval timers, along with the tool-specific action types.

The Collector already includes the ability to find the instrumentation points in a program, report them to the Client, and instantiate point actions at the Client’s request. It can also instantiate point actions for a set of instrumentation points that match a pattern chosen by the user or the tool builder (such as, “entry to all functions whose names start with `MPI_`”).

### 4.2.2 Static Instrumentation

Tools that rely on instrumentation that is compiled or linked directly into an application do not need to declare action types, but they do need to declare Client database attributes, and they also need to define at least one callback.

When the tool initializes itself in the Collector, it declares the attributes in the usual way, and then instead of defining action types and callbacks to be instantiated when data arrives from the program, it declares a single callback that the Collector will execute when the target program terminates. This routine can read a file that the program’s instrumentation has written and forward the data it collects to the Client using Tool Gear’s standard functions.

The instrumentation can choose a unique name for the file based on the program name and Unix process id. This information is also available to the Collector, so it can easily find the relevant file. For parallel programs, the instrumentation may write all the data to a single file, or it may write one file per process. The Collector can handle both situations. When any process in a parallel job terminates, DPCL invokes the callback and sends the process id as a parameter. Therefore, the Collector can read individual files as they are created. If only one file is created, it can be identified with the process id of task zero in the parallel job, and the callback can simply return when it is invoked for other processes in the job.

Tools that use dynamic instrumentation can easily determine what part of a program generated a piece of data. This allows the Collector to tell the Client database how to associate data with a program’s source code. Tools that do not use dynamic instrumentation must find another way to associate data with program locations. One method is for the instrumentation to decode the target program’s symbol table information so that it can look up file, function, and line number information based on the program counter. GNU libraries are available to help applications do this. The instrumentation then stores source location information in its output file along with the data it collects. When the callback function in the Collector reads this file, it must parse this information and declare the files, functions, and program locations to the Client database. It must also generate a unique tag for each program location. It will use this tag when sending data, and the Client database will then associate the data correctly with source code.

## 4.3 Creating the User Interface

The final step in defining a tool is to specify how the user will request actions and how data will be displayed.

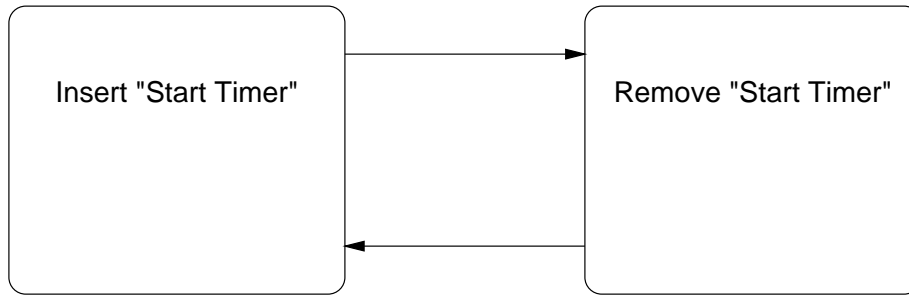


Figure 4: Tool developers define state diagrams for each action using the Tool Programming Interface. Along with the basic states and transitions, developers specify corresponding icons and menu text.

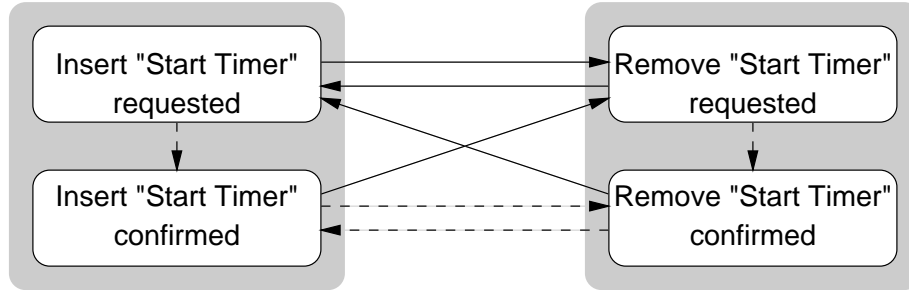


Figure 5: Tool Gear expands the basic state diagram that the tool developer declares to include transitional states. End users can initiate only those transitions indicated by solid arrows, while the Collector can normally initiate only the transitions indicated by dashed lines. In error conditions, the Client or Collector can force the system into any state.

#### 4.3.1 Requesting Actions

Tool builders specify how to request actions through the Tool Programming Interface. Our model for user interactions with target programs is that each tool defines one or more actions that the user can request, and each request is associated with a location in the target program. (The “program locations” in the Client typically represent DPCL instrumentation points. We use separate terms to distinguish the objects that DPCL defines and manipulates on the target computer from those that the Client uses to represent them.) Examples of actions include inserting (or removing) some instrumentation or setting a breakpoint.

When a tool runs, tool-specific software makes calls to the Tool Programming Interface to define a set of actions, along with associated icons, menu entries, and help text. The tool also defines how actions relate to each other by defining a state transition graph (see Figure 4.) The tool uses this graph for all the program locations (or a subset of them, depending on the tool), but the Client maintains separate states for each location. The current state for a program location determines what icon is displayed, what menu choices are available to the user, and what help text is displayed when the cursor passes over an icon. The tool also defines a default initial state.

Tool Gear can automatically expand the graph that the tool defines to differentiate between *requested* actions and *confirmed* actions. It also defines additional transitions between the new states (Figure 5.) Expanding the state diagram allows the display to show different icons for requested and confirmed actions, so the user can see when a pending action has been completed. As Figure 5 shows, users can only initiate requests, while the Collector normally only confirms them.

For example (again using Figure 5), when a tool starts up, each

program location will be in a default initial state of **Remove “Start Timer” confirmed**, meaning that no instrumentation has been inserted yet. From here, the only transition that the user can initiate for this action is to **Insert “Start Timer” requested**. Menu text that appears when the user clicks on a program location will reflect this choice, and choosing this text will cause a transition to that state. The Client will then send the corresponding request to the Collector update the icon displayed at the program location. From this new state, the user now has the choice of requesting removal of the (not yet confirmed) timer. Selecting this option would tell the Collector to cancel the pending request. The Collector, meanwhile, can change the state to **Insert “Start Timer” confirmed**. Each state transition causes the GUI to display the corresponding tool-defined icon.

Now suppose that the user does request cancellation of a pending action. The state transition would be from **Insert “Start Timer” requested** to **Remove “Start Timer” requested**. The Collector will receive the cancellation request (which will look like an ordinary request to remove the instrumentation). If the Collector is single-threaded, it may not be possible to cancel the action in progress. Therefore, the Collector will complete the insertion request and then immediately carry out the removal request. After completing the first action, it will send a confirmation back to the Client. At this point, the Client will be in the **Remove “Start Timer” requested** state, but the message will confirm an insertion. Since there is no transition allowed from the current state to to **Insert “Start Timer” confirmed**, the Client will ignore this message, and the icon on the display will not change. Then when the Client receives the removal confirmation, it will move to the **Remove “Start Timer” confirmed** state, and the icon appropri-



ate icon will be displayed. We designed this behavior specifically to avoid the potentially confusing situation that arises when a user makes a request and then changes his mind: if the GUI updated the icon after each confirmation arrived from the Collector, the display would show that the request was confirmed *after* the user had cancelled it. Then, a moment later, the display would show that the request had indeed been cancelled. We prefer to avoid this intermediate transition.

If the Collector cannot complete a request for some reason, it can send a negative acknowledgment message. Since this is an error condition, the Client can make a transition to any state. This allows it to show the true current state of the request.

The expanded transition diagram also includes Collector-initiated transitions between the confirmed states. When the user requests an action to be inserted on a set of program locations instead of one specific location, this transition allows the Collector to confirm the action at each location, even though the user didn't issue requests for the individual program locations.

A tool may support several independent groups of actions (such as inserting/removing instrumentation and setting/deleting breakpoints). The state diagrams that the tool defines for these groups of actions need not be connected. Thus, a program location can be in several independent states at the same time, reflecting the status of independent actions. The GUI automatically displays icons for independent states next to each other and builds the menus appropriately. For example, if the "Start Timer" action had been inserted at a particular location, clicking on that location might bring up a menu with the options of removing the "Start Timer" action, requesting a "Stop Timer" action, and requesting a breakpoint.

Although our model may seem complex, it needs to meet several needs:

- The user should be able to see the difference between a request and the outcome of a request. This corresponds to separating requested states from confirmed states.
- The user should not be able to force a program location into a confirmed state. This corresponds to defining separate transitions for the user and the Collector.
- The user should be able to cancel a pending request and to see the progress of the cancellation. This corresponds to allowing user transitions between two different requested states and omitting transitions from requested states to nonmatching confirmed states.
- The display should be able to respond to confirmations of multiple actions that were initiated by a single user request. This corresponds to defining transitions between confirmed states.

#### 4.3.2 Using the Database

The second task in defining a user interface is determining how data will be stored and displayed. When the Collector starts, tool-specific initialization code sends requests to the Client to define one or more columns in the database, as mentioned earlier. The Client then automatically creates and labels these columns in the Tree View. Later, when the Collector receives data from the target program, it forwards it to the Client with information describing the row and column where it should be stored in the database. New values for a given cell can either replace existing ones or be added to them.

For parallel programs, the database stores data for a given row and column separately for each thread or process. The Tree View display presents a single value for a cell, which can be (at the user's

choice) the sum, mean, minimum, maximum, standard deviation, or count of the values for that row and column. The tool can specify which of these summary values is displayed by default for each column. Placing the cursor over a value shows all the summary data in a line at the bottom of the window.

The Tree View further summarizes data in each column by displaying any one of the statistics mentioned above for the function, file, and program. This could show, for example, the total time for all functions that were instrumented in a program. Again, the user can select which type of summary to display. The Client computes all these values on the fly as data arrives in the database, so there is no delay when the user chooses different values to display.

At present, we have only implemented the Tree View display. However, we have begun to work on ways to present data graphically (in bar charts and pie charts, for example), and we expect to include that capability in future versions of Tool Gear.

## 5 Tools Built with Tool Gear

We have built two prototype tools using the Tool Gear infrastructure. The first, called TGmpx, displays cache utilization and FLOP data for selected sections of code (Figure 6.) It uses dynamic instrumentation. The second, called TGmpip, shows the cost of certain MPI calls. It uses instrumentation linked in through MPI's standard profiling interface.

### 5.1 TGmpx

TGmpx uses data gathered by the MPX hardware counter multiplexing library [9]. This library collects data from hardware performance counters, and on systems whose counter architecture would otherwise prevent concurrent counting of certain combinations of events, it uses time sharing to sample counters and extrapolate results. The MPX library also includes simple functions that can be compiled into an application to report the cache utilization and FLOP rate for specified section of code.

For the TGmpx tool, these functions were turned into DPCL probe modules. We defined two action types: one to start the measurement, and one to stop it and report the results. (We later added an interval timer to this tool; Section 5.3 describes that feature.) Starting the counters requires no callback function in the Collector. Stopping the counters invokes a callback in the Collector that simply forwards the data values (cache hit rate, FLOP count, and FLOP rate) back to the Client, along with a tag that identifies the instrumentation point where the callback was invoked and a process/thread identifier. The Client stores the data in the database and presents the results on the appropriate source code line.

For this tool, we defined a third action type to implement simple breakpoints. When the user sets a breakpoint, DPCL installs an instrumentation function that sends an empty message back to the Collector and then puts the process to sleep for a short time. When the callback in the Collector receives this message, it tells DPCL to pause the program's execution. The sleep call in the instrumentation function gives the Collector time to receive the message and call DPCL before the program continues to the next instruction.

The Client's Tree View display marks program locations with dotted boxes, and users can pop up a menu at any of these locations to install one of the three action types. A corresponding icon will then appear at that location. In Figure 6, the icon that appears before the function names indicates a request to start the counters, and the icon after the function names indicates a request to stop the counters and report results.

Figure 6 shows data for an eight-process parallel program. The

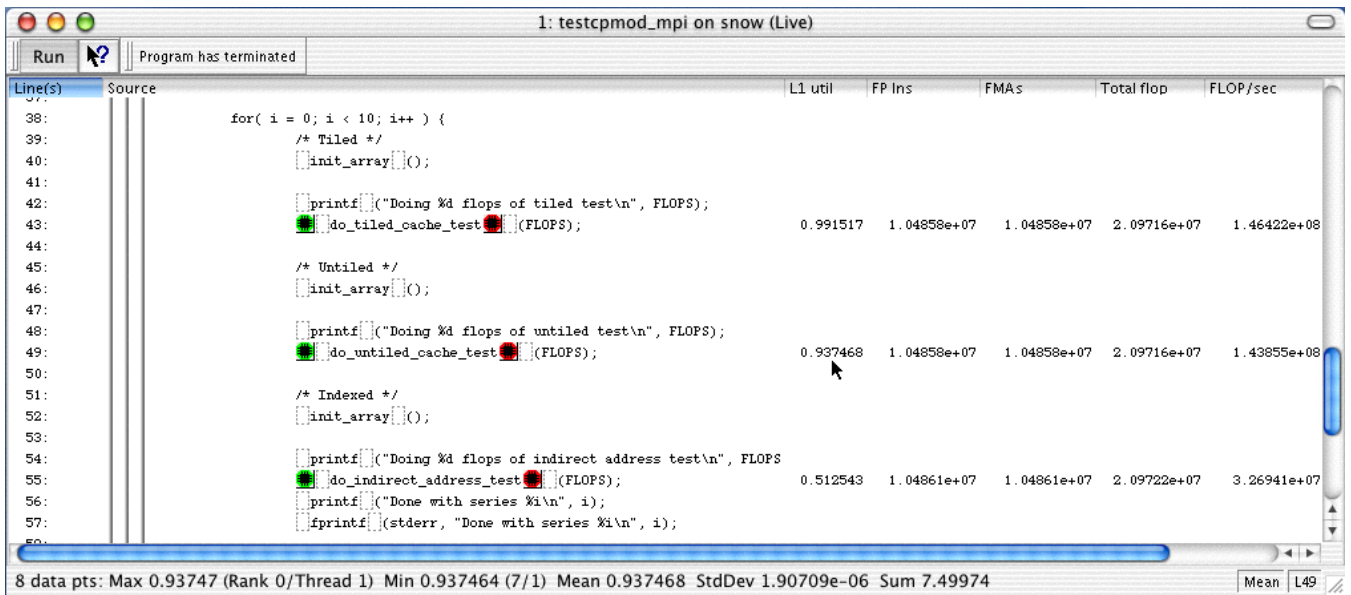


Figure 6: The TGmpx tool shows cache utilization and FLOP rates for selected regions of code in an eight-process parallel program. Here, three function calls have been instrumented (at runtime), and resulting data is shown on the corresponding source lines. A statistical summary of the cell under the cursor appears at the bottom of the display.

tool also works for sequential programs and multithreaded programs.

All of the display capability is built into the Tool Gear infrastructure. This includes not only the features described so far but also the ability to search in the source code, customizable tool tips that present help text for various elements of the display, and a programmable “About...” box. The parts of the user interface that are unique to TGmpx are the icons, the menu text (not shown here), and the column labels.

The Tree View’s collapsible display not only makes the source code for a large program easier to navigate, it also improves performance. The Client doesn’t ask the Collector for a list of functions in a file until the user clicks on the file name to open that display, nor does it request instrumentation points or source code until the user asks to see them. This approach greatly improves the scalability of the system, since it eliminates the delays (up to several minutes) that would occur if the Tree View requested all the information from a large executable at once.

## 5.2 TGmpip

TGmpip displays data generated by mpiP [17], a tool that instruments MPI calls and writes out a summary file describing how much time certain calls took to complete. It helps users find the communication calls in their codes that are taking a disproportionate amount of time.

This tool relies on linked-in instrumentation, so it does not need any action type declarations. Instead, we simply use DPCL to run the program. When the target finishes execution, the Collector invokes a callback to find and read the output file (whose name is based on the executable name, the number of processes, and the process id—all information that is available within the Collector). This file lists call sites for each instrumented MPI function, and for each function, there is a per-process listing of the number of times the function was called; the minimum, maximum, and mean execution time; and the percentage of time this MPI call accounts for in the total communication time and total application running time.

The tool-specific callback function culls this information from the file and transmits the per-call information to the Client for display (Figure 7).

Because this tool outputs no data until the program has finished running, there is no need for breakpoints or dynamic instrumentation. Therefore, we initially considered implementing a Collector for this tool that didn’t use DPCL. However, we quickly realized that the new Collector would still need to start the parallel program, detect when it had finished executing, identify the processes, and serve source code to the Client. All of this is certainly feasible, but we decided it would be simpler to take advantage of the functionality that already exists in DPCL and not write a new version of the Collector that didn’t use it. As a result, the tool-independent part of the TGmpip Collector is the same as what TGmpx uses. The disadvantage of this approach is that it limits the Collector to running on systems where DPCL is available, currently IBM and Linux.

An alternative implementation of this tool could use dynamic instrumentation to avoid the need for linking the target program with an instrumentation library. Such a tool could present a continually-updated display as the program ran. Much of its functionality could be implemented by dynamically inserting interval timers at each call to an MPI function and keeping a count of the calls.

## 5.3 Tool Development Time

One of our main goals for Tool Gear is to simplify the development process so that ideas can be transformed rapidly into working, usable tools. Although the ease of implementing a tool with Tool Gear is difficult to quantify, the time needed to build a tool is a reasonable measure of our success.

Because TGmpx was developed concurrently with Tool Gear itself, it is difficult to estimate how long it would have taken to build TGmpx on top of the existing Tool Gear infrastructure.

For TGmpip, we needed to extend the Tool Gear infrastructure somewhat to handle data from static instrumentation. This took

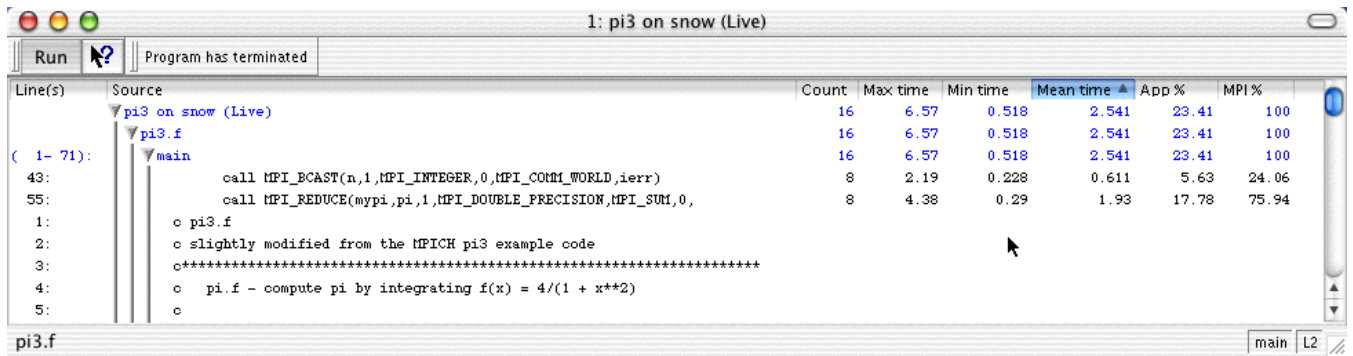


Figure 7: The TGmpip tool shows the cost of MPI communication calls. In this view, the user has sorted the display by Mean time column (highlighted), so the source lines with data in that column have been pulled out of the rest of the code and arranged in increasing order. The top three lines show the sum of the data in each column for the program, file, and function. Since the program consists of one function, all three are the same in this case. Unlike the example in Figure 6, the cursor is not pointing to a cell that contains data, so the line at the bottom of the window simply displays the file, function, and line where the cursor is pointing.

about a week, starting from the existing Tool Gear infrastructure and the stand-alone text-based tool that gathers the data.

Perhaps the best indication of how long it takes to create a new tool with Tool Gear is our experience with the interval timer tool. This tool simply measures and reports the elapsed time from a Start Timer call to a Stop Timer call. Users insert these calls using the standard dynamic instrumentation model we have described. A simplified version of the source code for the instrumentation appears in Figure 3. Implementing this tool took about half a day and required no changes to the infrastructure. The process closely followed the steps we have outlined in this paper.

We expect to further shorten the development time for tools by streamlining the programming interfaces as we gain more experience with them. We will also investigate putting all of the tool customization in the Collector, leaving a generic Client that tool builders would not have to customize unless they were building entirely new displays. In this scenario, the Collector's initialization function would describe to the Client not only actions and database attributes but also the state diagrams, icons, and help text. The Client already accepts a pathname for the Collector on its command line, so users can invoke different tools through the standard Client by calling customized scripts. The advantage of putting all the specialized code in the Collector is better encapsulation of the tool-specific functionality, which should simplify the development of tools.

#### 5.4 Future Tools

We have begun work on implementing new kinds of tools with Tool Gear. One of these is based on Umpire [16], which is a tool that automatically checks MPI programs for a variety of interprocess communication errors. This tool will display a list of suspected errors, and the user will be able to click on an error to get a display of the relevant source code. Although the error list will require a new display, much of the other needed functionality already exists in Tool Gear.

Another possible tool is a simple, lightweight memory leak detector. Using dynamic instrumentation to track calls to memory allocation and deallocation functions, this tool could keep track of the total number bytes allocated and freed. If the totals didn't match, the tool could show where the imbalance occurred. Of course, good memory analysis tools already exist, but they often require time-consuming instrumentation steps, and they can greatly increase ex-

ecution time. This new tool would excel at quick checks with minimal increases in run time.

## 6 Current Status and Future Work

Source code for Tool Gear and the sample tools TGmpx and TGmpip are currently available free of charge. (*A Web site will be given in the final paper.*) DPCL currently runs only on IBM and Linux systems, and the Linux port is experimental. However, IBM has released this code as open source, and the Dyninst technology on which it is based has been ported to many platforms. The Client, on the other hand, is already quite portable. It has been tested on IBM, Solaris, Linux, and Macintosh OS X workstations.

There are many areas for future improvements to Tool Gear, some of which we have mentioned elsewhere in this paper. As other developers begin to build tools with it, we will need to refine and extend the programming interfaces and add new functionality, especially for displaying data in different ways. We will also investigate extensions to the database. At a minimum, it should be able to write out data in a form that can easily be read into a spreadsheet program, so users can carry out more sophisticated data analyses. We will also consider adding some analysis capabilities to the database itself.

Our overall plan for Tool Gear is to discover what common features tool developers need, and then to extend Tool Gear so that developers can incorporate these features in their tools in a straightforward way. We will continue to take advantage of external open-source software packages as appropriate and to develop our own software components as needed.

With Tool Gear's flexibility, ease of use, and sophisticated interface features, we believe it has an exciting future as a foundation on which many useful new tools can be built.

## Acknowledgment

This work was performed under the auspices of the U.S. Department of Energy by University of California Lawrence Livermore National Laboratory under contract number W-7405-Eng-48. UCRL-JC-147901-EXT-ABS.

## References

- [1] DEPARTMENT OF COMPUTER AND INFORMATION SCIENCE, UNIVERSITY OF OREGON. *TAU User's Guide*. Eugene, Oregon, 2000. <http://www.acl.lanl.gov/tau>.
- [2] DEROSE, L., HOOVER JR., T., AND HOLLINGSWORTH, J. K. The Dynamic Probe Class Library—An infrastructure for developing instrumentation for performance tools. In *Proceedings 15th International Parallel and Distributed Processing Symposium* (April 2001).
- [3] DEROSE, L., AND REED, D. A. SvPablo: A multi-language architecture-independent performance analysis system. In *Proceedings of the International Conference on Parallel Processing* (September 1999).
- [4] DEROSE, L. A., ZHANG, Y., AYDT, R., PANTANO, M., AND WHITMORE, S. *SvPablo User's Guide*. Department of Computer Science, University of Illinois, Urbana, Illinois, November 2001. <http://www-pablo.cs.uiuc.edu>.
- [5] GYLLENHAAL, J. C., HWU, W. W., AND RAU, B. R. HMDES version 2 specification. Tech. Rep. IMPACT-96-03, University of Illinois, Urbana, Illinois, 1996. <http://www.crhc.uiuc.edu/~gyllen/>.
- [6] HOLLINGSWORTH, J. K., MILLER, B. P., AND CARGILLE, J. Dynamic program instrumentation for scalable performance tools. In *Proceedings of the Scalable High Performance Computing Conference* (May 1994), pp. 841–850.
- [7] KOHL, J. A., AND CASAVANT, T. Use of PARADISE: A meta-tool for visualizing parallel systems. In *Proceedings of the Fifth International Parallel Processing Symposium* (1991), pp. 561–567.
- [8] MAY, J., AND BERMAN, F. Retargetability and extensibility in a parallel debugger. *Journal of Parallel and Distributed Computing* 35, 2 (June 1996), 142–155.
- [9] MAY, J. M. MPX: Software for multiplexing hardware performance counters in multithreaded programs. In *Proceedings 15th International Parallel and Distributed Processing Symposium* (April 2001).
- [10] MILLER, B. P., CALLAGHAN, M. D., CARGILLE, J. M., HOLLINGSWORTH, J. K., IRVIN, R. B., KARAVANIC, K. L., KUNCHITHAPADAM, K., AND NEWHALL, T. The Paradyn parallel performance measurement tools. *IEEE Computer* 28, 11 (November 1995), 37–46.
- [11] MOHR, B., BROWN, D., AND MALONEY, A. TAU: A portable parallel program analysis environment for pC++. In *Proceedings of CONPAR94—VAPP VI* (September 1994), pp. 29–40.
- [12] PARADYN PROJECT. *Paradyn Parallel Performance Tools User's Guide*. Computer Science Department, University of Wisconsin, Madison, Wisconsin, January 2002. <http://www.cs.wisc.edu/paradyn/>.
- [13] SOCHA, D., BAILEY, M. L., AND NOTKIN, D. Voyeur: Graphical views of parallel programs. *Proceedings of the ACM SIGPLAN/SIGOPS Workshop on Parallel and Distributed Debugging, published in ACM SIGPLAN NOTICES* 24, 1 (January 1989), 206–215.
- [14] STASKO, J. T., AND KRAEMER, E. A methodology for building application-specific visualizations of parallel programs. *Journal of Parallel and Distributed Computing* 18, 2 (June 1993), 258–264.
- [15] TROLLTECH. *Qt Reference Documentation*. <http://doc.trolltech.com>.
- [16] VETTER, J. S., AND DE SUPINSKI, B. R. Dynamic software testing of MPI applications with Umpire. In *Proceedings SC2000* (November 2000).
- [17] VETTER, J. S., AND MCCracken, M. O. Statistical scalability analysis of communication operations in distributed applications. In *Eighth ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming* (June 2001).